**Google DeepMind**

# Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities.

Gemini Team, Google

In this report, we introduce the Gemini 2.X model family: Gemini 2.5 Pro and Gemini 2.5 Flash, as well as our earlier Gemini 2.0 Flash and Flash-Lite models. Gemini 2.5 Pro is our most capable model yet, achieving SoTA performance on frontier coding and reasoning benchmarks. In addition to its incredible coding and reasoning skills, Gemini 2.5 Pro is a thinking model that excels at multimodal understanding and it is now able to process up to 3 hours of video content. Its unique combination of long context, multimodal and reasoning capabilities can be combined to unlock new agentic workflows. Gemini 2.5 Flash provides excellent reasoning abilities at a fraction of the compute and latency requirements and Gemini 2.0 Flash and Flash-Lite provide high performance at low latency and cost. Taken together, the Gemini 2.X model generation spans the full Pareto frontier of model capability vs cost, allowing users to explore the boundaries of what is possible with complex agentic problem solving.

## 1. Introduction

We present our latest family of natively multimodal models with advanced reasoning through thinking, long context and tool-use capabilities: Gemini 2.5 Pro and 2.5 Flash and our earlier Gemini 2.0 Flash and Gemini 2.0 Flash-Lite models. Together these form a new family of highly-capable models representing our next generation of AI models, designed to power a new era of agentic systems. Building upon the foundation of the Gemini 1.5 series (Gemini Team, 2024), this Gemini 2.X generation brings us closer to the vision of a universal AI assistant (Hassabis, 2025).

The Gemini 2.X series are all built to be natively multimodal, supporting long context inputs of >1 million tokens and have native tool use support. This allows them to comprehend vast datasets and handle complex problems from different information sources, including text, audio, images, video and even entire code repositories. These extensive capabilities can also be combined to build complex agentic systems, as happened in the case of Gemini Plays Pokémon[1] (Zhang, 2025). Different models in the series have different strengths and capabilities: (1) Gemini 2.5 Pro is our most intelligent thinking model, exhibiting strong reasoning and code capabilities. It excels at producing interactive web applications, is capable of codebase-level understanding and also exhibits emergent multimodal coding abilities. (2) Gemini 2.5 Flash is our hybrid reasoning model with a controllable thinking budget, and is useful for most complex tasks while also controlling the tradeoff between quality, cost, and latency. (3) Gemini 2.0 Flash is our fast and cost-efficient non-thinking model for everyday tasks and (4) Gemini 2.0 Flash-Lite is our fastest and most cost-efficient model, built for at-scale usage. A full comparison of the models in the Gemini 2.X model family is provided in Table 1. Taken together, the Gemini 2.X family of models cover the whole Pareto frontier of model capability vs cost, shifting it forward across a large variety of core capabilities, applications and use-cases, see Figure 1.

The Gemini 2.5 family of models maintain robust safety metrics while improving dramatically on

---

[1]Pokémon is a trademark of Nintendo Co., Ltd., Creatures Inc., and Game Freak Inc.

---

| | *Gemini 1.5 Flash* | *Gemini 1.5 Pro* | **Gemini 2.0 Flash-Lite** | **Gemini 2.0 Flash** | **Gemini 2.5 Flash** | **Gemini 2.5 Pro** |
|---|---|---|---|---|---|---|
| **Input modalities** | Text, Image, Video, Audio | Text, Image, Video, Audio | Text, Image, Video, Audio | Text, Image, Video, Audio | Text, Image, Video, Audio | Text, Image, Video, Audio |
| **Input length** | 1M | 2M | 1M | 1M | 1M | 1M |
| **Output modalities** | Text | Text | Text | Text, Image* | Text, Audio* | Text, Audio* |
| **Output length** | 8K | 8K | 8K | 8K | 64K | 64K |
| **Thinking** | No | No | No | Yes* | Dynamic | Dynamic |
| **Supports tool use?** | No | No | No | Yes | Yes | Yes |
| **Knowledge cutoff** | November 2023 | November 2023 | June 2024 | June 2024 | January 2025 | January 2025 |

Table 1 | Comparison of Gemini 2.X model family with Gemini 1.5 Pro and Flash. Tool use refers to the ability of the model to recognize and execute function calls (e.g., to perform web search, complete a math problem, execute code). *currently limited to Experimental or Preview, see Section 2.7. Information accurate as of publication date.*

helpfulness and general tone compared to their 2.0 and 1.5 counterparts. In practice, this means that the 2.5 models are substantially better at providing safe responses without interfering with important use cases or lecturing end users. We also evaluated Gemini 2.5 Pro's Critical Capabilities, including CBRN, cybersecurity, machine learning R&D, and deceptive alignment. While Gemini 2.5 Pro showed a significant increase in some capabilities compared to previous Gemini models, it did not reach any of the Critical Capability Levels in any area.

Our report is structured as follows: we begin by briefly describing advances we have made in model architecture, training and serving since the release of the Gemini 1.5 model. We then showcase the performance of the Gemini 2.5 models, including qualitative demonstrations of its abilities. We conclude by discussing the safety evaluations and implications of this model series.

## 2. Model Architecture, Training and Dataset

### 2.1. Model Architecture

The Gemini 2.5 models are sparse mixture-of-experts (MoE) (Clark et al., 2022; Du et al., 2021; Fedus et al., 2021; Jiang et al., 2024; Lepikhin et al., 2020; Riquelme et al., 2021; Roller et al., 2021; Shazeer et al., 2017) transformers (Vaswani et al., 2017) with native multimodal support for text, vision, and audio inputs. Sparse MoE models activate a subset of model parameters per input token by learning to dynamically route tokens to a subset of parameters (experts); this allows them to decouple total model capacity from computation and serving cost per token. Developments to the model architecture contribute to the significantly improved performance of Gemini 2.5 compared to Gemini 1.5 Pro (see Section 3). Despite their overwhelming success, large transformers and sparse MoE models are known to suffer from training instabilities (Chowdhery et al., 2022; Dehghani et al., 2023; Fedus et al., 2021; Lepikhin et al., 2020; Liu et al., 2020; Molybog et al., 2023; Wortsman et al., 2023; Zhai et al., 2023; Zhang et al., 2022). The Gemini 2.5 model series makes considerable progress in enhancing large-scale training stability, signal propagation and optimization dynamics, resulting in a considerable boost in performance straight out of pre-training compared to previous Gemini models.
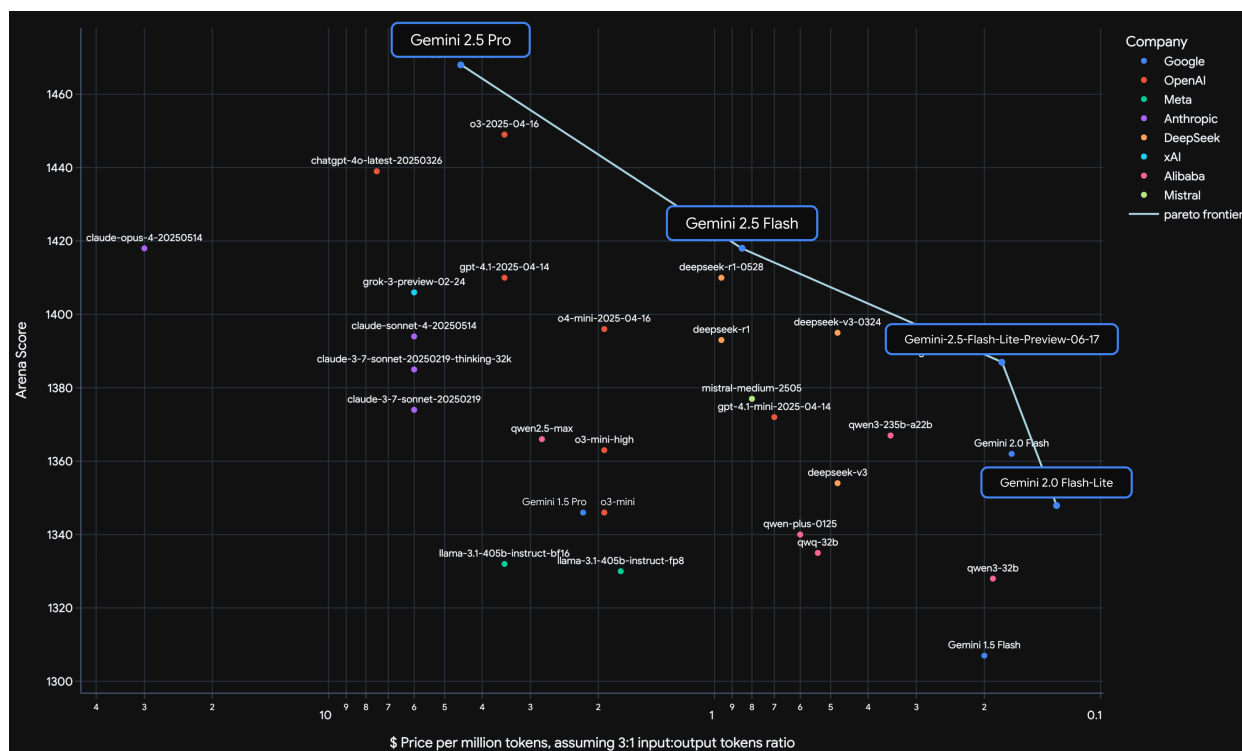
Figure 1 | Cost-performance plot. Gemini 2.5 Pro is a marked improvement over Gemini 1.5 Pro, and has an LMArena score that is over 120 points higher than Gemini 1.5 Pro. Cost is a weighted average of input and output tokens pricing per million tokens. Source: LMArena, imported on 2025-06-16.

Gemini 2.5 models build on the success of Gemini 1.5 in processing long-context queries, and incorporate new modeling advances allowing Gemini 2.5 Pro to surpass the performance of Gemini 1.5 Pro in processing long context input sequences of up to 1M tokens (see Table 3). Both Gemini 2.5 Pro and Gemini 2.5 Flash can process pieces of long-form text (such as the entirety of "Moby Dick" or "Don Quixote"), whole codebases, and long form audio and video data (see Appendix 8.5). Together with advancements in long-context abilities, architectural changes to Gemini 2.5 vision processing lead to a considerable improvement in image and video understanding capabilities, including being able to process 3-hour-long videos and the ability to convert demonstrative videos into interactive coding applications (see our recent blog post by Baddepudi et al., 2025).

The smaller models in the Gemini 2.5 series — Flash size and below — use distillation (Anil et al., 2018; Hinton et al., 2015), as was done in the Gemini 1.5 series (Gemini Team, 2024). To reduce the cost associated with storing the teacher's next token prediction distribution, we approximate it using a k-sparse distribution over the vocabulary. While this still increases training data throughput and storage demands by a factor of k, we find this to be a worthwhile trade-off given the significant quality improvement distillation has on our smaller models, leading to high-quality models with a reduced serving cost (see Figure 2).

## 2.2. Dataset

Our pre-training dataset is a large-scale, diverse collection of data encompassing a wide range of domains and modalities, which includes publicly available web documents, code (various programming languages), images, audio (including speech and other audio types) and video, with a cutoff date of June 2024 for 2.0 and January 2025 for 2.5. Compared to the Gemini 1.5 pre-training dataset